# Genome Hacking

**Yaniv Erlich
Whitehead Institute for
Biomedical Research**

**Twitter: @erlichya**

# Public data is important for genetic studies

**Research** ─────────────────

GENOME RESEARCH

## Exome sequencing and disease-network analysis of a single family implicate a mutation in *KIF1A* in hereditary spastic paraparesis

Yaniv Erlich,[1,4,5] Simon Edvardson,[2,4] Emily Hodges,[3] Shamir Zenvirt,[2] Pramod Thekkat,[3] Avraham Shaag,[2] Talya Dor,[2] Gregory J. Hannon,[3] and Orly Elpeleg[2]

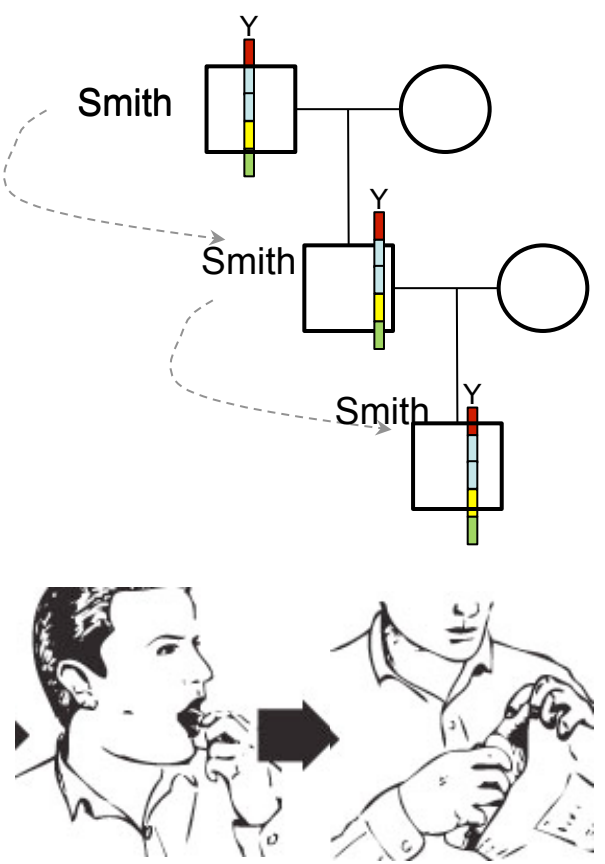**AJHG** The American Journal of Human Genetics

**REPORT** ───────

## Joubert Syndrome 2 (JBTS2) in Ashkenazi Jews Is Associated with a *TMEM216* Mutation

Simon Edvardson,[1,9] Avraham Shaag,[2,9] Shamir Zenvirt,[3] Yaniv Erlich,[5,6] Gregory J. Hannon,[5,6] Alan L. Shanske,[8] John Moshe Gomori,[4] Joseph Ekstein,[7] and Orly Elpeleg[2,3,*]

**To make this endeavor sustainable, we must proactively map risks**

# Co-segregation between Y-chr and surnames



**www.ysearch.org:**

# Exploiting genetic genealogy databases

## The Washington Post

### Found on the Web, With DNA: a Boy's Father

By Rob Stein
Washington Post Staff Writer
Sunday, November 13, 2005

Like many children whose mothers used an anonymous sperm donor, the 15-year-old boy longed for any shred of information about his biological father. But, uniquely, this resourceful teenager decided to try exploiting the latest in genetic technology and the sleuthing powers of the Internet in his quest.

By submitting a DNA sample to a commercial genetic database service designed to help people draw their family tree, the youth found a crucial clue that quickly enabled him to track down his long-sought parent.

"I was stunned," said Wendy Kramer, whose online registry for children trying to find anonymous donors of sperm or egg helped lead the teenager to his father. "This had never been done before. No one knew you could get a DNA test and find your donor."

While welcomed by advocates of children trying to locate anonymous donors, the case -- apparently the first of its kind -- has raised alarm among sperm banks and some medical ethicists. They are concerned it might start a trend that could violate the privacy of thousands of sperm donors and discourage future ones.

An anecdote?

# The main idea – a <span style="color:red">systematic</span> study

**Can we recover the identity of anonymous sequencing datasets using public resources?**

# Empirical test: what is the probability to recover a surname?



**Y-STR of a real person**

**Querying Ysearch and SMGF**

**Calculating surname confidence score**

**Inferring surname**

**x900**

**Comparing the predicted surname to the true one**

Expectation for US Caucasian males from middle and upper class:
12% Successful recoveries

# The Venter case

lobSTR

- **We got a surname from whole genome sequencing data**

  **Method:** lobSTR: A short tandem repeat profiler for personal genomes

  Melissa Gymrek,[1,2] David Golan,[2,3] Saharon Rosset,[3] and Yaniv Erlich[2,4]

  [1]Harvard-MIT Division of Health Sciences and Technology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA; [2]Whitehead Institute for Biomedical Research, Cambridge, Massachusetts 02142, USA; [3]Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 69978, Israel

- **The DNA does not belong to Craig Venter**

GENOME RESEARCH

DYS 458
17

**Try it yourself: bit.ly/craig_venter_haplotype_updated**

# Can we identify anonymous personal genomes?

# Recovering the identifies of CEU individuals

1000 Genomes
A Deep Catalog of Human Genetic Variation

EMBL-EBI

**10 CEU genomes**

SORENSON MOLECULAR GENEALOGY FOUNDATION

y*search*

8 Surname predictions with Utah ancestry

+You   **Search**   Images   Maps   Play   YouTube   News   Gmail   Documents   Calendar   More -

Google        Winfield Utah

Found an obituary that has the exact description of the pedigree

**CORIELL INSTITUTE**
FOR MEDICAL RESEARCH

**Probability of a random match < 5x10$^{-9}$**

# Beginner's luck?



p<5x10$^{-9}$    p<5x10$^{-6}$    p<10$^{-5}$

■ Successful surname recovery (targeted individual)

↗ Person tested by genetic genealogy service (source)

— Patrilineal line from source to target

Breaching the privacy of close to **50** CEU samples.

# Summary

Our approach:

- No experimental work involved.
- The identifying information propagates via deep genealogical ties.
- The attack completely relies on public resources.

Testing close to 1000 Y-STR haplotypes, demonstrating complete identification of Venter and close to 50 CEU individuals.

# IMHO, recommendations

1. Consent:
      - Be honest about risks. Be honest about benefits.

2. Multi-tier approach:
      -  Give participants options for data sharing.

3. Proactive approach:
      - Keep mapping risks. Friendly hacking is far better than a real one.

4.   Technical solutions:
      - We did not explore those enough. Much more to do here.

# Acknowledgements



**Melissa Gymrek** (HST – Harvard/MIT)
Amy McGuire (Baylor)
David Golan (Tel-Aviv University)
Eran Halperin (Tel-Aviv University)



Science
AAAS

## Identifying Personal Genomes by Surname Inference

Melissa Gymrek,[1,2,3,4] Amy L. McGuire,[5] David Golan,[6] Eran Halperin,[7,8,9] Yaniv Erlich[1]*

Sharing sequencing data sets without identifiers has become a common practice in genomics. Here, we report that surnames can be recovered from personal genomes by profiling short tandem repeats on the Y chromosome (Y-STRs) and querying recreational genetic genealogy databases. We show that a combination of a surname with other types of metadata, such as age and state, can be used to triangulate the identity of the target. A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources. We quantitatively analyze the probability of identification for U.S. males. We further demonstrate the feasibility of this technique by tracing back with high probability the identities of multiple participants in public sequencing projects.

Open Access (with FREE registration)

**Funding:**
**Andria and Paul Heafy**
**Jim and Cathy Stone**